

```
In [1]: import pandas as pd
import numpy as np

data = pd.read_csv("Bengaluru_House_Data.csv")

In [2]: data.head()
```

	area_type	availability	location	size	society	total_sqft	bath	balcony	price
0	Super built-up Area	Ready To Move	Electronic City Phase II	2 BHK	Coomee	1056	2.0	1.0	39.07
1	Plot Area	Ready To Move	Chikka Tripathi	4 Bedroom	Theanp	2600	5.0	3.0	120.00
2	Built-up Area	Ready To Move	Uttarahalli	3 BHK	NaN	1440	2.0	3.0	62.00
3	Super built-up Area	Ready To Move	Lingadheeranahalli	3 BHK	Sowere	1521	3.0	1.0	95.00
4	Super built-up Area	Ready To Move	Kothanur	2 BHK	NaN	1200	2.0	1.0	51.00

In [4]: data.shape

Out[4]: (13320, 9)

In [5]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13320 entries, 0 to 13319
Data columns (total 9 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   area_type   13320 non-null    object
 1   availability 13320 non-null    object
 2   location    13319 non-null    object
 3   size        13304 non-null    object
 4   society     7818 non-null    object
 5   total_sqft  13320 non-null    float64
 6   bath        13247 non-null    float64
 7   balcony     12921 non-null    float64
 8   price       13320 non-null    float64
dtypes: float64(3), object(6)
memory usage: 624.4+ KB
```

In [6]: #Value counts for each columns

```
for column in data.columns:
    print(data[column].value_counts())
    print("=====")

Super built-up Area    8799
Built-up Area          2438
Plot Area              2025
Carpet Area            87
Name: area_type, dtype: int64
=====
Ready To Move          18581
18-Dec                 307
18-May                 295
18-Apr                 271
18-Aug                 266
...
15-Jun                 1
16-Jan                 1
16-Oct                 1
16-Nov                 1
14-Jul                 1
Name: availability, Length: 81, dtype: int64
=====
Whitefield             540
Sarjapur Road          399
Electronic City         382
Kansapura Road         273
Thanisandra            234
...
Near ulias theater     1
Near Electronic City,  1
kg halli jaihalli west 1
sarjapura main road    1
Jagadish Nagar         1
Name: location, Length: 1305, dtype: int64
=====
2 BHK                  5199
3 BHK                  4310
4 Bedroom              826
4 BHK                  591
3 Bedroom              547
1 BHK                  538
2 Bedroom              329
5 Bedroom              297
6 Bedroom              191
1 Bedroom              185
8 Bedroom              84
7 Bedroom              83
5 BHK                  59
9 Bedroom              46
6 BHK                  38
7 BHK                  17
1 RK                   13
10 Bedroom             12
9 BHK                  8
8 BHK                  5
10 BHK                 2
11 BHK                 2
11 Bedroom             1
13 BHK                 1
27 BHK                 1
14 BHK                 1
19 BHK                 1
16 BHK                 1
12 Bedroom             1
18 Bedroom             1
43 Bedroom             1
Name: size, dtype: int64
=====
GrrvaGr               89
PrarePa               76
Sryalan               59
Prlates               59
Grown E               56
...
Vbhonj                1
Fae Zov               1
Adace P               1
VarLisa               1
Nehshee               1
Name: society, Length: 2688, dtype: int64
=====
1280                   843
1180                   221
1580                   285
2480                   196
680                    189
...
3124                   ...
42889                  1
5665.84                1
673                    1
4480 - 6640            1
Name: total_sqft, Length: 2117, dtype: int64
=====
2.0                    6988
3.0                    3286
4.0                    1226
1.0                    788
5.0                    524
6.0                    273
7.0                    182
8.0                    64
9.0                    43
10.0                   13
12.0                   7
13.0                   3
11.0                   3
16.0                   2
27.0                   1
18.0                   1
40.0                   1
15.0                   1
14.0                   1
Name: bath, dtype: int64
=====
2.0                    5119
1.0                    4897
3.0                    1672
6.0                    1029
Name: balcony, dtype: int64
=====
75.00                  310
65.00                  302
55.00                  278
60.00                  270
46.00                  148
...
81.55                   1
69.49                   1
42.18                   1
70.25                   1
74.82                   1
Name: price, Length: 1994, dtype: int64
=====
```

In [7]: data.isnull().sum()

	area_type	availability	location	size	society	total_sqft	bath	balcony	price
0	Super built-up Area	Ready To Move	Electronic City Phase II	2 BHK	Coomee	1056	2.0	1.0	39.07
1	Plot Area	Ready To Move	Chikka Tripathi	4 Bedroom	Theanp	2600	5.0	3.0	120.00
2	Built-up Area	Ready To Move	Uttarahalli	3 BHK	NaN	1440	2.0	3.0	62.00
3	Super built-up Area	Ready To Move	Lingadheeranahalli	3 BHK	Sowere	1521	3.0	1.0	95.00
4	Super built-up Area	Ready To Move	Kothanur	2 BHK	NaN	1200	2.0	1.0	51.00

In [8]: data.drop(columns = ['area\_type','availability','society','balcony'], inplace = True)

In [9]: data.describe()

	bath	price
count	13247.000000	12320.000000
mean	2.892610	112.565627
std	1.341458	148.971674
min	1.000000	8.000000
25%	2.000000	50.000000
50%	2.000000	72.000000
75%	3.000000	120.000000
max	40.000000	3600.000000

In [10]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13320 entries, 0 to 13319
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   location    13319 non-null    object
 1   size        13304 non-null    object
 2   total_sqft  13320 non-null    float64
 3   bath        13247 non-null    float64
 4   price       13320 non-null    float64
dtypes: float64(2), object(3)
memory usage: 364.3+ KB
```

In [11]: data['location'].value\_counts()

	bath	price
Whitefield	540	540
Sarjapur Road	399	399
Electronic City	382	382
Kansapura Road	273	273
Thanisandra	234	234
...	...	...
Near ulias theater	1	1
Near Electronic City,	1	1
kg halli jaihalli west	1	1
sarjapura main road	1	1
Jagadish Nagar	1	1
Name: location, Length: 1305, dtype: int64		

In [12]: #There is only one missing value in location so we are filling it with Sarjapur Road

```
data['location'] = data['location'].fillna('Sarjapur Road')
```

In [13]: data['size'].value\_counts()

	bath	price
2 BHK	5199	5199
3 BHK	4310	4310
4 Bedroom	826	826
4 BHK	591	591
3 Bedroom	547	547
1 BHK	538	538
2 Bedroom	329	329
5 Bedroom	297	297
6 Bedroom	191	191
1 Bedroom	185	185
8 Bedroom	84	84
7 Bedroom	83	83
5 BHK	59	59
9 Bedroom	46	46
6 BHK	38	38
7 BHK	17	17
1 RK	13	13
10 Bedroom	12	12
9 BHK	8	8
8 BHK	5	5
10 BHK	2	2
11 BHK	2	2
11 Bedroom	1	1
13 BHK	1	1
27 BHK	1	1
14 BHK	1	1
19 BHK	1	1
16 BHK	1	1
12 Bedroom	1	1
18 Bedroom	1	1
43 Bedroom	1	1
Name: size, dtype: int64		

In [14]: #There are 16 null values in size hence we are filling those values with 2BHK

```
data['size'] = data['size'].fillna("2 BHK")
```

In [15]: #We have 73 missing values in bathroom so we are placing values with median

```
data['bath'] = data['bath'].fillna(data['bath'].median())
```

In [16]: data.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 13320 entries, 0 to 13319
Data columns (total 5 columns):
 #   Column      Non-Null Count  Dtype
---  --
 0   location    13320 non-null    object
 1   size        13320 non-null    object
 2   total_sqft  13320 non-null    float64
 3   bath        13320 non-null    float64
 4   price       13320 non-null    float64
dtypes: float64(2), object(3)
memory usage: 364.3+ KB
```

In [17]: data['bhk'] = data['size'].str.split().str.get(0).astype(int)

In [18]: #data.bhk > 20

```
#There are 2 location where we have more than 27 BHK, So these are outliers in data so we need to fix it.
```

	location	size	total_sqft	bath	price	bhk
1718	Electronic City Phase II	27 BHK	8000	27.0	230.0	27
4684	Munnekolal	43 Bedroom	2400	40.0	660.0	43

In [19]: data['total\_sqft'].unique()

```
#We have got the values in form of ranges we need to fix it so we will take mean of ot like(1133+1384)/2
```

Out[19]: array(['1656', '2689', '1440', ..., '1133 + 1384', '774', '4689'], dtype=object)

In [20]: def convertRange(x):

```
temp = x.split(",")
if len(temp) == 2:
    return (float(temp[0]) + float(temp[1]))/2
try:
    return float(x)
except:
    return None
```

In [21]: data['total\_sqft'] = data['total\_sqft'].apply(convertRange)

In [22]: data.head()

	location	size	total_sqft	bath	price	bhk
0	Electronic City Phase II	2 BHK	1056.0	2.0	39.07	2
1	Chikka Tripathi	4 Bedroom	2600.0	5.0	120.00	4
2	Uttarahalli	3 BHK	1440.0	2.0	62.00	3
3	Lingadheeranahalli	3 BHK	1521.0	3.0	95.00	3
4	Kothanur	2 BHK	1200.0	2.0	51.00	2

## price per square feet

In [23]: #Price/total\_sqft=price per square ft

```
data['price_per_sqft'] = data['price']/100000/data['total_sqft']
```

In [24]: data['price\_per\_sqft']

	total_sqft	bath	price	bhk	price_per_sqft
count	13073.000000	13320.000000	13320.000000	13320.000000	1.307300e+04
mean	1554.942029	2.688814	112.565627	2.802778	7.949600e+03
std	1238.458773	1.338754	148.971674	1.294496	1.072400e+05
min	1.000000	1.000000	8.000000	1.000000	2.670298e+02
25%	1100.000000	2.000000	50.000000	2.000000	4.265734e+03
50%	1275.000000	2.000000	72.000000	3.000000	5.454545e+03
75%	1670.000000	3.000000	120.000000	3.000000	7.338057e+03
max	52272.000000	40.000000	3600.000000	43.000000	1.200000e+07

In [25]: data['location'].value\_counts()

	total_sqft	bath	price	bhk	price_per_sqft
count	13073.000000	13320.000000	13320.000000	13320.000000	1.307300e+04
mean	1554.942029	2.688814	112.565627	2.802778	7.949600e+03
std	1238.458773	1.338754	148.971674	1.294496	1.072400e+05
min	1.000000	1.000000	8.000000	1.000000	2.670298e+02
25%	1100.000000	2.000000	50.000000	2.000000	4.265734e+03
50%	1275.000000	2.000000	72.000000	3.000000	5.454545e+03
75%	1670.000000	3.000000	120.000000	3.000000	7.338057e+03
max	52272.000000	40.000000	3600.000000	43.000000	1.200000e+07

Out[25]:

	total_sqft	bath	price	bhk	price_per_sqft
count	13073.000000	13320.000000	13320.000000	13320.000000	1.307300e+04
mean	1554.942029	2.688814	112.565627	2.802778	7.949600e+03
std	1238.458773	1.338754	148.971674	1.294496	1.072400e+05
min	1.000000	1.000000	8.000000	1.000000	2.670298e+02
25%	1100.000000	2.000000	50.000000	2.000000	4.265734e+03
50%	1275.000000	2.000000	72.000000	3.000000	5.454545e+03
75%	1670.000000	3.000000	120.000000	3.000000	7.338057e+03
max	52272.000000	40.000000	3600.000000	43.000000	1.200000e+07

In [26]: data['location'].value\_counts()

	total_sqft	bath	price	bhk	price_per_sqft
Whitefield	540	540	540	540	540
Sarjapur Road	399	399	399	399	399
Electronic City	382	382	382	382	382
Kansapura Road	273	273	273	273	273
Thanisandra	234	234	234	234	234
...	...	...	...	...	...
Chokkasandra	1	1	1	1	1
Rahmath Nagar	1	1	1	1	1
Whitefield	1	1	1	1	1
Sonam Layout	1	1	1	1	1
Escorts Colony	1	1	1	1	1
Name: location, Length: 1306, dtype: int64					

In [27]: #Any location having < 10 value\_counts replace the location with others

```
data['location'] = data['location'].apply(lambda x: x.strip())
location_count = data['location'].value_counts()
```

In [28]: location\_count\_less\_10 = location\_count[location\_count<=10]

```
#This are the location having <10 counts
```

Out[28]:

	total_sqft	bath	price	bhk	price_per_sqft
Sector 1 HSR Layout	10	10	10	10	10
1st Block Koramangala	10	10	10	10	10
Sadashiva Nagar	10	10	10	10	10
Nagandhapura	10	10	10	10	10
Nagadevanahalli	10	10	10	10	10
...	...	...	...	...	...
Chokkasandra	1	1	1	1	1
Rahmath Nagar	1	1	1	1	1
Sonam Layout	1	1	1	1	1
Swaraa Nagar	1	1	1	1	1
Escorts Colony	1	1	1	1	1
Name: location, Length: 1064, dtype: int64					

In [29]: data['location'] = data['location'].apply(lambda x: 'other' if x in location\_count\_less\_10 else x)

In [30]: data['location'].value\_counts()

```
#So we got 2886 other values which was having <10 vlaue_counts
```

Out[30]:

	total_sqft	bath	price	bhk	price_per_sqft
other	2886	2886	2886	2886	2886
Whitefield	540	540	540	540	540
Sarjapur Road	399	399	399	399	399
Electronic City	382	382	382	382	382
Kansapura Road	273	273	273	273	273
...	...	...	...	...	...
Tindru	11	11	11	11	11
Nehru Nagar	11	11	11	11	11
Pattandur Agrahara	11	11	11	11	11
2nd Phase Judicial Layout	11	11	11	11	11
Thyagaraja Nagar	11	11	11	11	11
Name: location, Length: 242, dtype: int64					

## Outlier Detection and removal

In [31]: data.describe()

	total_sqft	bath	price	bhk	price_per_sqft
count	13073.000000	13320.000000	13320.000000	13320.000000	1.307300e+04
mean	1554.942029	2.688814	112.565627	2.802778	7.949600e+03
std	1238.458773	1.338754	148.971674	1.294496	1.072400e+05
min	1.000000	1.000000	8.000000	1.000000	2.670298e+02
25%	1100.000000	2.000000	50.000000	2.000000	4.265734e+03
50%	1275.000000	2.000000	72.000000	3.000000	5.454545e+03
75%	1670.000000	3.000000	120.000000	3.000000	7.338057e+03
max	52272.000000	40.000000	3600.000000	43.000000	1.200000e+07

In [32]: #data whose values price/bhk <300 that should not be an feeseble flat

```
data = data[(data['total_sqft']/data['bhk'])>=300]
```

Out[32]:

	total_sqft	bath	price	bhk	price_per_sqft
count	12329.000000	12329.000000	12329.000000	12329.000000	1.232900e+04
mean	1590.166773	2.561441	111.444236	2.651472	6322.476758
std	1261.827604	1.072551	152.759322	0.973754	4187.479096
min	300.000000	1.000000	8.440000	1.000000	267.829813
25%	1110.000000	2.000000	49.340000	2.000000	4207.119741
50%	1300.000000	2.000000	70.000000	3.000000	5300.000000
75%	1700.000000	3.000000	115.000000	3.000000	6938.463540
max	52272.000000	16.000000	3600.000000	16.000000	176470.588235

In [33]: data.shape

Out[33]: (12329, 7)

In [34]: data.price\_per\_sqft.describe()

	total_sqft	bath	price	bhk	price_per_sqft
count	12329.000000	12329.000000	12329.000000	12329.000000	1.232900e+04
mean	1590.166773	2.561441	111.444236	2.651472	6322.476758
std	1261.827604	1.072551	152.759322	0.973754	4187.479096
min	300.000000	1.000000	8.440000	1.000000	267.829813
25%	1110.000000	2.000000	49.340000	2.000000	4207.119741
50%	1300.000000	2.000000	70.000000	3.000000	5300.000000
75%	1700.000000	3.000000	115.000000	3.000000	6938.463540
max	52272.000000	16.000000	3600.000000	16.000000	176470.588235

<